

# SVRE: NEW METHOD FOR TRAINING GANs

GAUTHIER GIDEL

Mila, Université de Montréal  
Research intern at Element AI

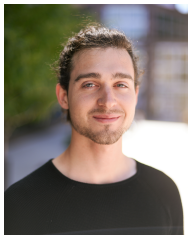
**GENERATIVE MODELING AND MODEL-BASED REASONING  
FOR ROBOTICS AND AI WORKSHOP**

June 14, 2019

# REDUCING NOISE IN GAN TRAINING WITH VARIANCE REDUCED EXTRAGRADIENT



TATJANA CHAVDAROVA \*



GAUTHIER GIDEL \*



FRANÇOIS FLEURET



SIMON LACOSTE-JULIEN

# GENERATIVE ADVERSARIAL NETWORKS

[Goodfellow et al., 2014]

# CHALLENGES

- Standard supervised learning:

$$\min_{\theta} \mathcal{L}(\theta)$$

- GANs: Hard (different) optimization problem: minimax.

# CHALLENGES

- Standard supervised learning:

$$\min_{\theta} \mathcal{L}(\theta)$$

- GANs: Hard (different) optimization problem: minimax.

$$\min_{\theta_G} \max_{\theta_D} V(\theta_G, \theta_D)$$

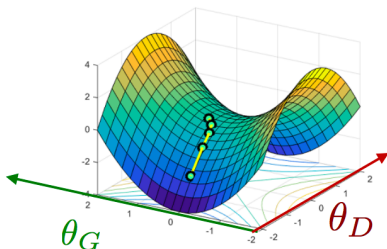
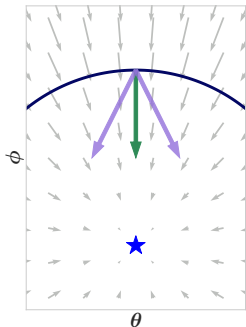


Image source: Vaishnavh Nagarajan

# “NOISE”: NOISY GRADIENT ESTIMATES DUE TO STOCHASTICITY

- Using sub-samples (mini-batches) of the full dataset to update the parameters
- Variance Reduced (VR) Gradient: optimization methods that reduce such noise

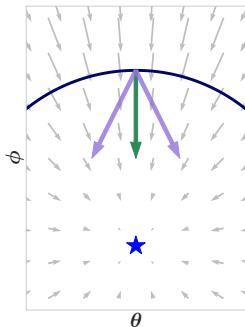
Minimization: Single-objective



- Batch method direction
- Stochastic method direction: noisy

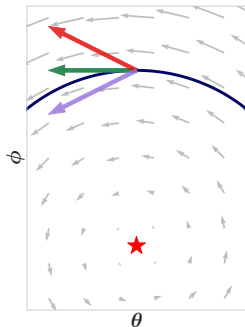
# VARIANCE REDUCTION—MOTIVATION FOR GAMES

- INTUITIVELY: **MINIMIZATION** VS. **GAME** (NOISE FROM STOCHASTIC GRADIENT)
- EMPIRICALLY: **BIGGAN**—“INCREASED BATCH SIZE SIGNIFICANTLY IMPROVES PERFORMANCES”
- TO SUM UP, TWO ISSUES:



Minimization

“approximately” the right direction



Game

Direction with noise can be “bad”

# VARIANCE REDUCTION–MOTIVATION FOR GAMES

- INTUITIVELY: **MINIMIZATION** VS. **GAME** (NOISE FROM STOCHASTIC GRADIENT)
- EMPIRICALLY: **BIGGAN**–“INCREASED BATCH SIZE SIGNIFICANTLY IMPROVES PERFORMANCES”
- TO SUM UP, TWO ISSUES:

Brock et al. [2018] report a relative improvement of **46%** of the Inception Score metric [Salimans et al., 2016] on **ImageNet** if the mini-batch size is increased **8**-fold.



# VARIANCE REDUCTION–MOTIVATION FOR GAMES

- INTUITIVELY: **MINIMIZATION** VS. **GAME** (NOISE FROM STOCHASTIC GRADIENT)
- EMPIRICALLY: **BIGGAN**–“INCREASED BATCH SIZE SIGNIFICANTLY IMPROVES PERFORMANCES”
- TO SUM UP, TWO ISSUES:
  - Adversarial aspect from min-max → **Extragradient**.
  - Noise from stochastic gradient → **Variance Reduction**.

# EXTRAGRADIENT

# EXTRAGRADIENT

Two players  $\theta, \varphi$ . Idea: perform a “Lookahead step”

$$\text{Extrapolation: } \begin{cases} \theta_{t+1/2} = \theta_t - \eta \nabla_{\theta} \mathcal{L}_G(\theta_t, \varphi_t) \\ \varphi_{t+1/2} = \varphi_t - \eta \nabla_{\varphi} \mathcal{L}_D(\theta_t, \varphi_t) \end{cases}$$

$$\text{Update: } \begin{cases} \theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}_G(\theta_{t+1/2}, \varphi_{t+1/2}) \\ \varphi_{t+1} = \varphi_t - \eta \nabla_{\varphi} \mathcal{L}_D(\theta_{t+1/2}, \varphi_{t+1/2}) \end{cases}$$

# VARIANCE REDUCED GRADIENT METHODS

# VARIANCE REDUCED ESTIMATE OF THE GRADIENT

Based on Finite sum assumption:

$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, \omega),$$

Epoch based algorithm:

- Save the full gradient  $\frac{1}{n} \sum_i \nabla \mathcal{L}(\mathbf{x}_i, \omega^S)$  and the snapshot  $\omega^S$ .
- For one epoch use the update rule:

$$\omega \leftarrow \omega - \eta \left[ \underbrace{\nabla \mathcal{L}(\mathbf{x}_i, \omega)}_{\text{Stochastic gradient}} + \underbrace{\frac{1}{n} \sum_i \nabla \mathcal{L}(\mathbf{x}_i, \omega^S) - \nabla \mathcal{L}(\mathbf{x}_i, \omega^S)}_{\text{correction using saved past iterate}} \right]$$

- Requires 2 stochastic gradients (at the current point and at the snapshot).
- If  $\omega^S$  is close to  $\omega \rightarrow$  close to full batch gradient  $\rightarrow$  small variance.
- Full batch gradient expensive but tractable, e.g., compute it once per pass.

# VARIANCE REDUCED ESTIMATE OF THE GRADIENT

Based on Finite sum assumption:

$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, \omega),$$

Epoch based algorithm:

- Save the full gradient  $\frac{1}{n} \sum_i \nabla \mathcal{L}(\mathbf{x}_i, \omega^S)$  and the snapshot  $\omega^S$ .
- For one epoch use the update rule:

$$\omega \leftarrow \omega - \eta \left[ \underbrace{\nabla \mathcal{L}(\mathbf{x}_i, \omega)}_{\text{Stochastic gradient}} + \underbrace{\frac{1}{n} \sum_i \nabla \mathcal{L}(\mathbf{x}_i, \omega^S) - \nabla \mathcal{L}(\mathbf{x}_i, \omega^S)}_{\text{correction using saved past iterate}} \right]$$

- Requires 2 stochastic gradients (at the current point and at the snapshot).
- If  $\omega^S$  is close to  $\omega \rightarrow$  close to full batch gradient  $\rightarrow$  small variance.
- Full batch gradient expensive but tractable, e.g., compute it once per pass.

# VARIANCE REDUCED ESTIMATE OF THE GRADIENT

Based on Finite sum assumption:

$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, \omega),$$

Epoch based algorithm:

- Save the full gradient  $\frac{1}{n} \sum_i \nabla \mathcal{L}(\mathbf{x}_i, \omega^S)$  and the snapshot  $\omega^S$ .
- For one epoch use the update rule:

$$\omega \leftarrow \omega - \eta \left[ \underbrace{\nabla \mathcal{L}(\mathbf{x}_i, \omega)}_{\text{Stochastic gradient}} + \underbrace{\frac{1}{n} \sum_i \nabla \mathcal{L}(\mathbf{x}_i, \omega^S) - \nabla \mathcal{L}(\mathbf{x}_i, \omega^S)}_{\text{correction using saved past iterate}} \right]$$

- Requires 2 stochastic gradients (at the current point and at the snapshot).
- If  $\omega^S$  is close to  $\omega \rightarrow$  close to full batch gradient  $\rightarrow$  small variance.
- Full batch gradient expensive but tractable, e.g., compute it once per pass.

# VARIANCE REDUCED ESTIMATE OF THE GRADIENT

Based on Finite sum assumption:

$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, \omega),$$

Epoch based algorithm:

- Save the full gradient  $\frac{1}{n} \sum_i \nabla \mathcal{L}(\mathbf{x}_i, \omega^S)$  and the snapshot  $\omega^S$ .
- For one epoch use the update rule:

$$\omega \leftarrow \omega - \eta \left[ \underbrace{\nabla \mathcal{L}(\mathbf{x}_i, \omega)}_{\text{Stochastic gradient}} + \underbrace{\frac{1}{n} \sum_i \nabla \mathcal{L}(\mathbf{x}_i, \omega^S) - \nabla \mathcal{L}(\mathbf{x}_i, \omega^S)}_{\text{correction using saved past iterate}} \right]$$

- Requires 2 stochastic gradients (at the current point and at the snapshot).
- If  $\omega^S$  is close to  $\omega \rightarrow$  close to full batch gradient  $\rightarrow$  small variance.
- Full batch gradient expensive but tractable, e.g., compute it once per pass.



# SVRE: VARIANCE REDUCTION + EXTRAGRADIENT PSEUDO-ALGORITHM

1. Save snapshot  $\omega^S \leftarrow \omega_t$  and compute  $\frac{1}{n} \sum_i \nabla \mathcal{L}(\mathbf{x}_i, \omega^S)$ .
2. For  $i$  in  $1, \dots, \text{epoch\_length}$ :
  - Compute  $\omega_{t+\frac{1}{2}}$  with variance reduced gradients at  $\omega_t$ .
  - Compute  $\omega_{t+1}$  with variance reduced gradients at  $\omega_{t+\frac{1}{2}}$ .
  - $t \leftarrow t + 1$
3. Repeat until convergence.

# SVRE: VARIANCE REDUCTION + EXTRAGRADIENT PSEUDO-ALGORITHM

1. Save snapshot  $\omega^S \leftarrow \omega_t$  and compute  $\frac{1}{n} \sum_i \nabla \mathcal{L}(x_i, \omega^S)$ .
2. For  $i$  in  $1, \dots, \text{epoch\_length}$ :
  - Compute  $\omega_{t+\frac{1}{2}}$  with variance reduced gradients at  $\omega_t$ .
  - Compute  $\omega_{t+1}$  with variance reduced gradients at  $\omega_{t+\frac{1}{2}}$ .
  - $t \leftarrow t + 1$
3. Repeat until convergence.

# SVRE: VARIANCE REDUCTION + EXTRAGRADIENT PSEUDO-ALGORITHM

1. Save snapshot  $\omega^S \leftarrow \omega_t$  and compute  $\frac{1}{n} \sum_i \nabla \mathcal{L}(x_i, \omega^S)$ .
2. For  $i$  in  $1, \dots, \text{epoch\_length}$ :
  - Compute  $\omega_{t+\frac{1}{2}}$  with variance reduced gradients at  $\omega_t$ .
  - Compute  $\omega_{t+1}$  with variance reduced gradients at  $\omega_{t+\frac{1}{2}}$ .
  - $t \leftarrow t + 1$
3. Repeat until convergence.

# SVRE: VARIANCE REDUCTION + EXTRAGRADIENT PSEUDO-ALGORITHM

1. Save snapshot  $\omega^S \leftarrow \omega_t$  and compute  $\frac{1}{n} \sum_i \nabla \mathcal{L}(x_i, \omega^S)$ .
2. For  $i$  in  $1, \dots, \text{epoch\_length}$ :
  - Compute  $\omega_{t+\frac{1}{2}}$  with variance reduced gradients at  $\omega_t$ .
  - Compute  $\omega_{t+1}$  with variance reduced gradients at  $\omega_{t+\frac{1}{2}}$ .
  - $t \leftarrow t + 1$
3. Repeat until convergence.

# SVRE: VARIANCE REDUCTION + EXTRAGRADIENT PSEUDO-ALGORITHM

1. Save snapshot  $\omega^S \leftarrow \omega_t$  and compute  $\frac{1}{n} \sum_i \nabla \mathcal{L}(x_i, \omega^S)$ .
2. For  $i$  in  $1, \dots, \text{epoch\_length}$ :
  - Compute  $\omega_{t+\frac{1}{2}}$  with variance reduced gradients at  $\omega_t$ .
  - Compute  $\omega_{t+1}$  with variance reduced gradients at  $\omega_{t+\frac{1}{2}}$ .
  - $t \leftarrow t + 1$
3. Repeat until convergence.

# SVRE: VARIANCE REDUCTION + EXTRAGRADIENT PSEUDO-ALGORITHM

1. Save snapshot  $\omega^S \leftarrow \omega_t$  and compute  $\frac{1}{n} \sum_i \nabla \mathcal{L}(x_i, \omega^S)$ .
2. For  $i$  in  $1, \dots, \text{epoch\_length}$ :
  - Compute  $\omega_{t+\frac{1}{2}}$  with variance reduced gradients at  $\omega_t$ .
  - Compute  $\omega_{t+1}$  with variance reduced gradients at  $\omega_{t+\frac{1}{2}}$ .
  - $t \leftarrow t + 1$
3. Repeat until convergence.

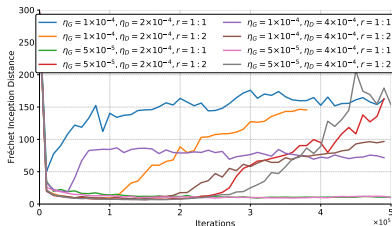
SVRE yields the fastest convergence rate for strongly convex stochastic game optimization in the literature.

## SVRE: EXPERIMENTS

# EXPERIMENTS

## SVRE YIELDS STABLE GAN OPTIMIZATION

### Stochastic baseline



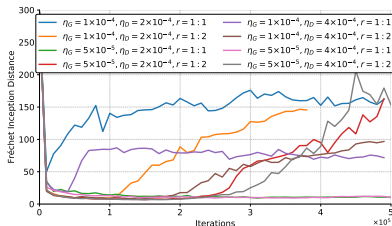
- Always diverges.
- Many hyperparameters  
( $\eta_G, \eta_D, \beta_1, \gamma, r$ ).
- + if convergence  $\rightarrow$  fast



# EXPERIMENTS

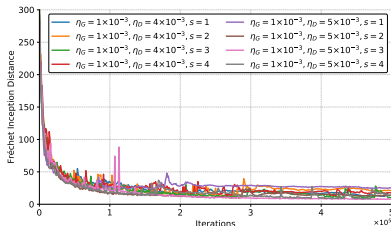
## SVRE YIELDS STABLE GAN OPTIMIZATION

### Stochastic baseline



- Always diverges.
- Many hyperparameters ( $\eta_G, \eta_D, \beta_1, \gamma, r$ ).
- + if convergence  $\rightarrow$  fast

### SVRE



- + Does not diverge.
- + fewer hyperparameters (omits  $\beta_1, \gamma, r$ )
- slower for very deep nets.

## SVRE: TAKEAWAYS

# SVRE: TAKEAWAYS

- Controlling variance is more critical for games (could be reason behind success of *Adam* on GANs)
- SVRE: combines Extragradient and variance reduction.
- Best convergence rate (under some assumptions) for games.
- Good stability properties.

THANKS.  
Questions?

# REFERENCES I

- A. Brock, J. Donahue, and K. Simonyan. Large Scale GAN Training for High Fidelity Natural Image Synthesis. ArXiv e-prints, September 2018.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Advances in Neural Information Processing Systems 27, pages 2672–2680. 2014.
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training GANs. In NIPS, 2016.

# APPENDIX

# THE GAN FRAMEWORK

EQUILIBRIUM AT  $p_g = p_d$

The discriminator maximizes:

$$\begin{aligned} V(G, D) &= \int_x p_d(x) \log(D(x)) dx + \int_z p_z(z) \log(1 - D(G(z))) dz \\ &= \int_x p_d(x) \log(D(x)) + p_g(x) \log(1 - D(x)) dx \end{aligned}$$

Where we used  $x = G(z)$ , and  $p_g$  is the distribution of  $x$ .  
Hence, the optimal discriminator  $D^*$  is:

$$D^*(x) = \frac{p_d(x)}{p_d(x) + p_g(x)}$$

# THE GAN FRAMEWORK

EQUILIBRIUM AT  $p_g = p_d$

The generator minimizes:

$$\begin{aligned}V(G, D^*) &= \mathbb{E}_{x \sim p_d} [\log D^*(x)] + \mathbb{E}_{x \sim p_g} [\log(1 - D^*(x))] \\&= \mathbb{E}_{x \sim p_d} \left[ \log \frac{p_d(x)}{p_d(x) + p_g(x)} \right] + \mathbb{E}_{x \sim p_g} \left[ \log \frac{p_g(x)}{p_d(x) + p_g(x)} \right] \\&= -\log 4 + \mathbb{D}_{KL}(p_d \parallel \frac{p_d + p_g}{2}) + \mathbb{D}_{KL}(p_g \parallel \frac{p_d + p_g}{2}) \\&= -\log 4 + 2 \cdot \mathbb{D}_{JS}(p_d \parallel p_g)\end{aligned}$$

where we used:  $\mathbb{D}_{JS}(p \parallel q) = \frac{1}{2} \mathbb{D}_{KL}(p \parallel \frac{p+q}{2}) + \frac{1}{2} \mathbb{D}_{KL}(q \parallel \frac{p+q}{2})$ .

The optimum is reached when  $p_g = p_d$ ,  
and the optimal value is  $-\log 4$ .